

# 美团四层负载均衡 - MGW

王伟@云计算部 20161008



# 自我介绍

---



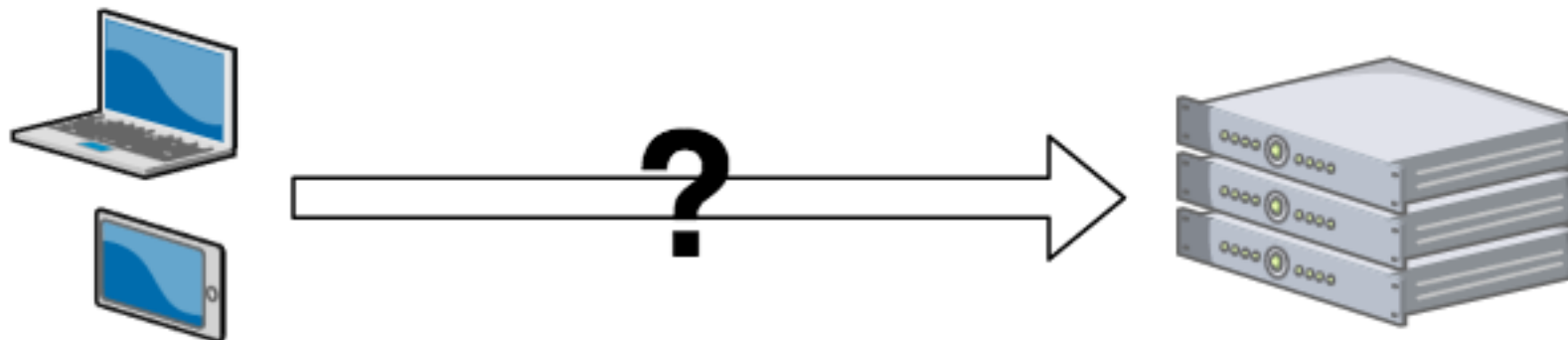
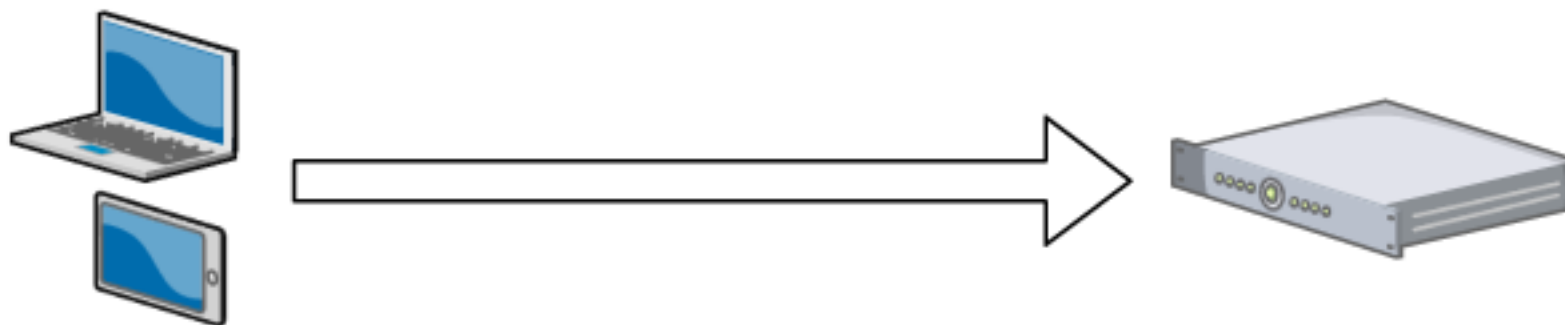
# 目录

---

- 负载均衡介绍
- 高性能
- 高可靠
- 技术展望

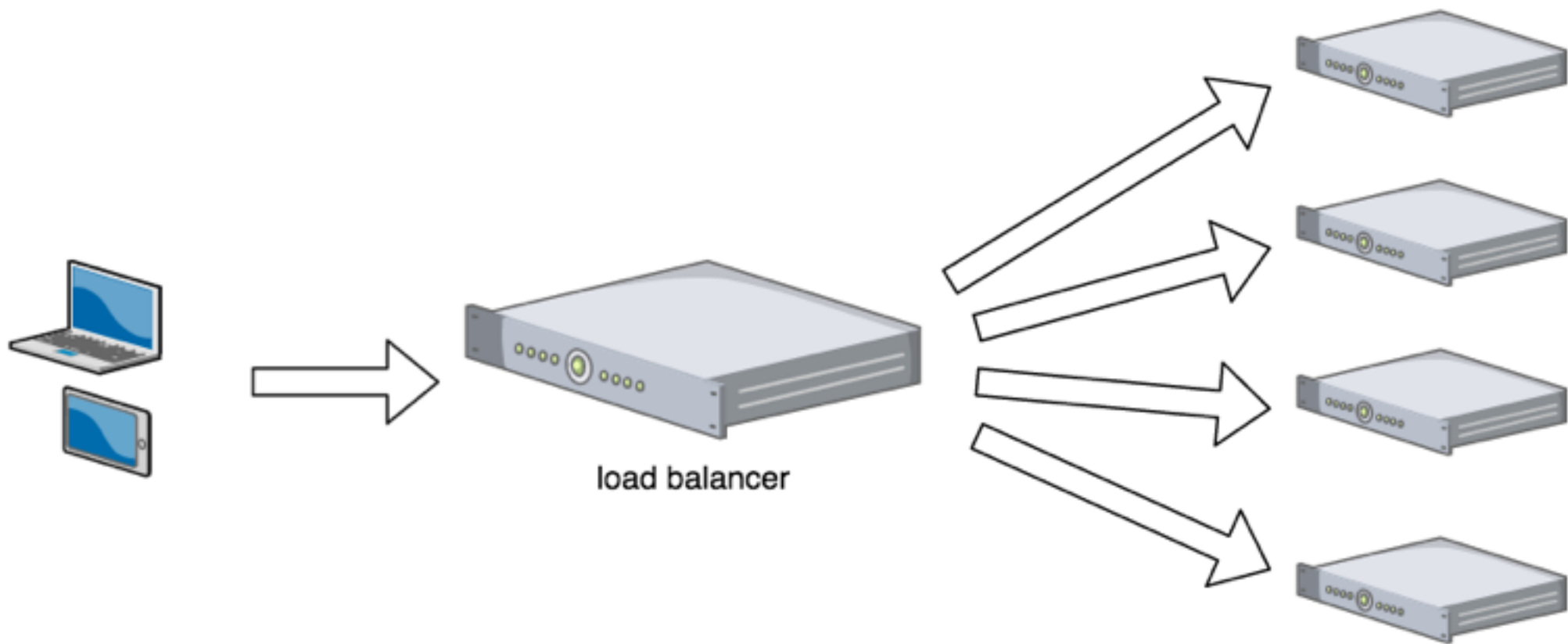
# 什么是负载均衡?

---



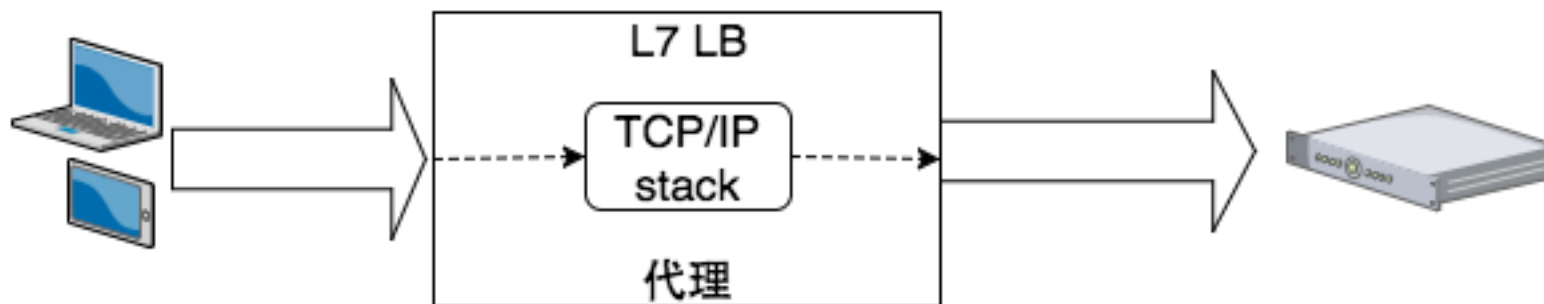
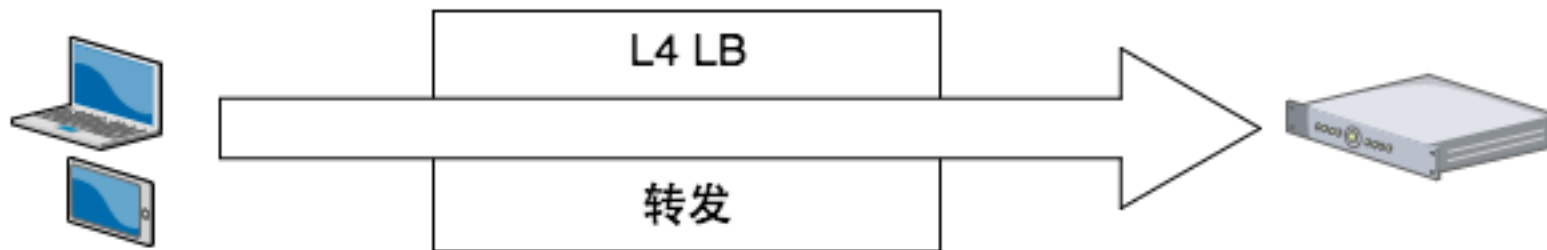
# 什么是负载均衡?

---



# 四层负载均衡 and 七层负载均衡?

---

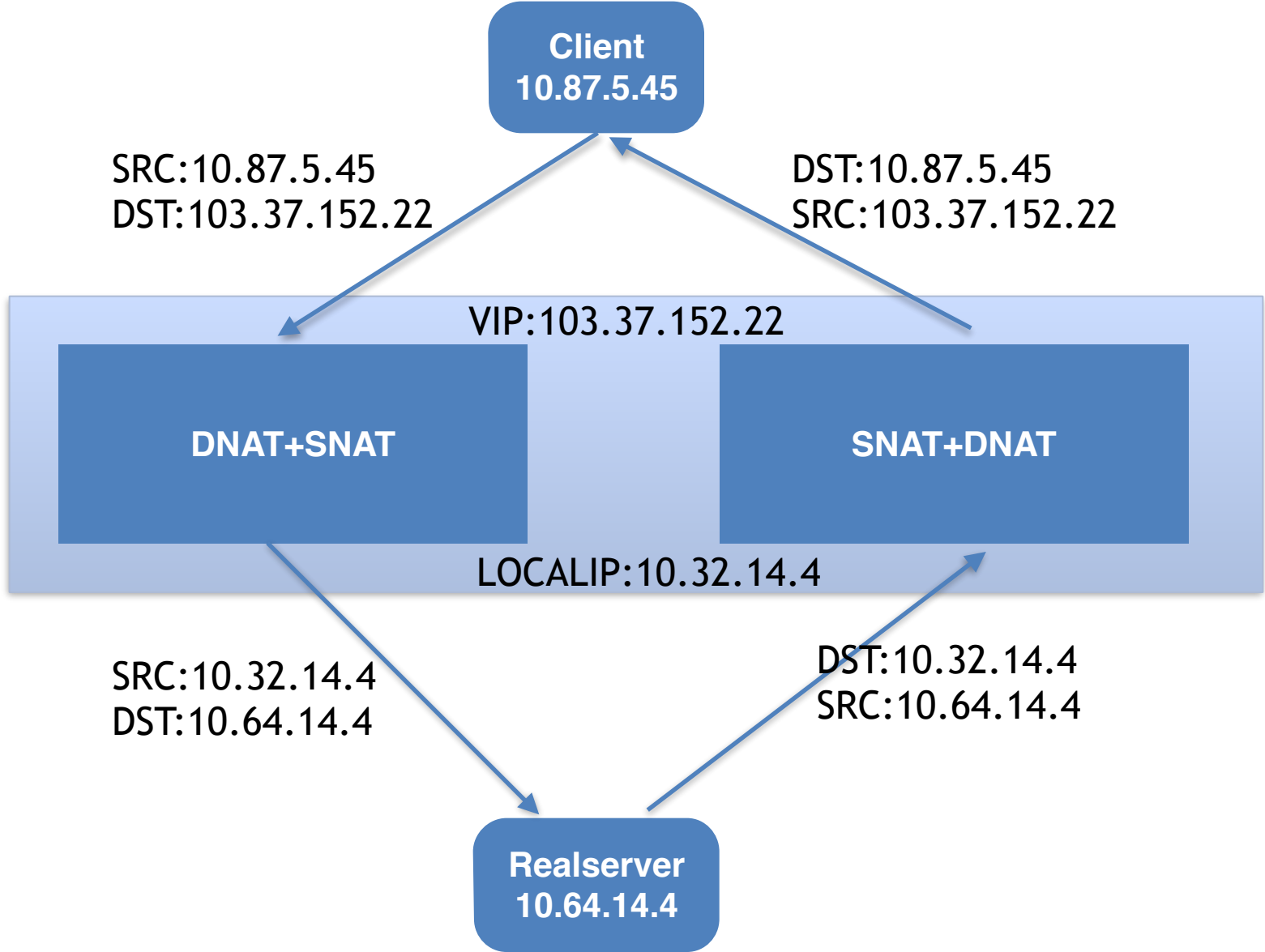


# 转发模式对比

模式	优点	缺点
DR(三角传输)	1. 应用直接将应答发给客户端，性能好	1. 必须在一个二层 2. 应用服务器需要配置VIP
NAT (DNAT)	1. 应用服务器无需做配置	1. 负载均衡必须以网关形式存在
TUNNEL	1. 和DR一样，应用直接将应答发给客户端，性能好。	1. 对应用服务器要求高，需要支持tunnel 2. 应用服务器需要配置vip
FULLNAT (SNAT+DNAT)	1. 应用服务器无需做配置 2. 对网络环境要求比较低	1. 丢失client ip



# 转发模式 - FULLNAT





# 目录

---

- 负载均衡介绍
- 高性能
- 高可靠
- 技术展望

# 负载均衡

.....

硬件负载均衡	软件负载均衡
	 LVS
	
	

# 硬件负载均衡

---

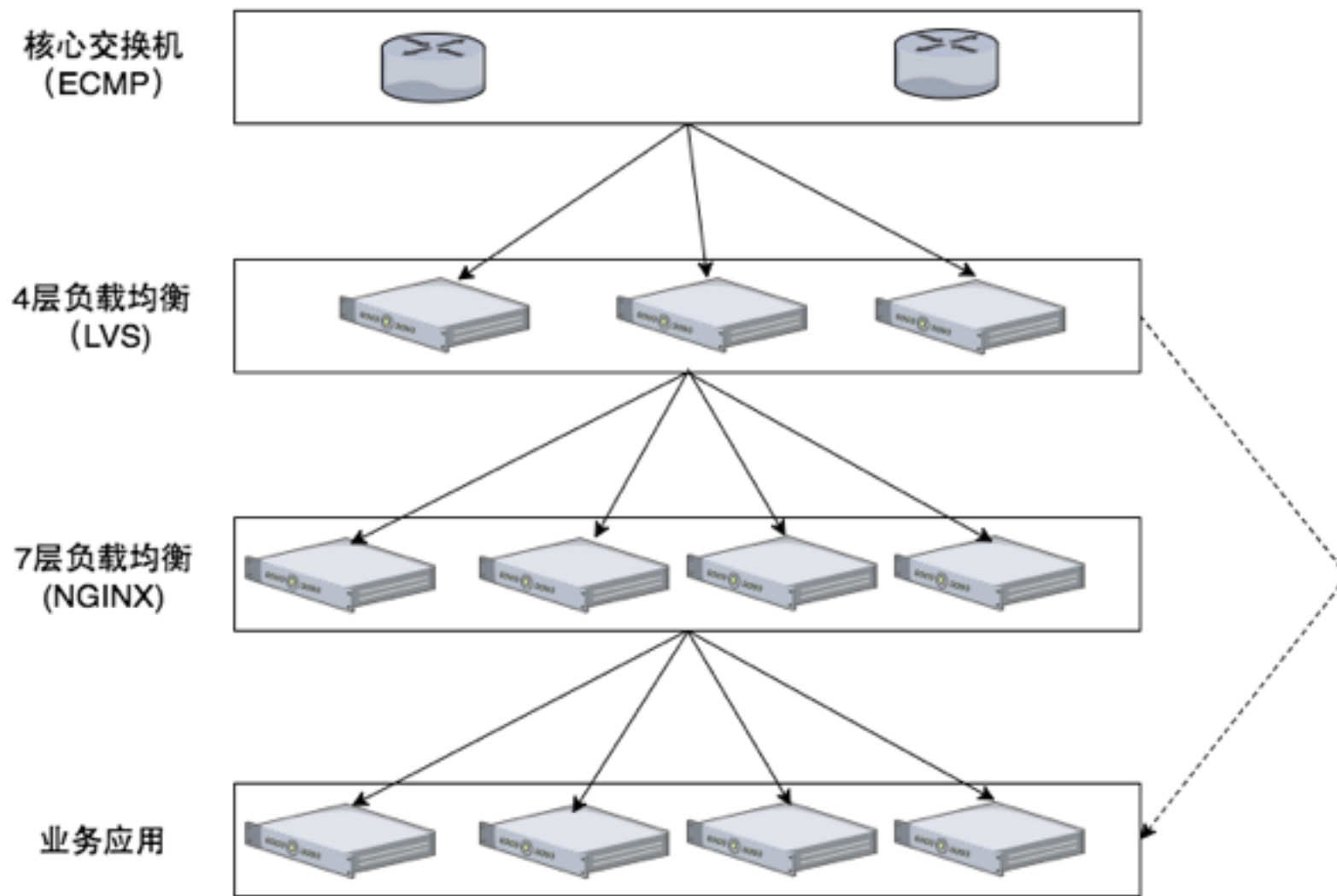
**1. 硬件成本**

**2. 人力成本**

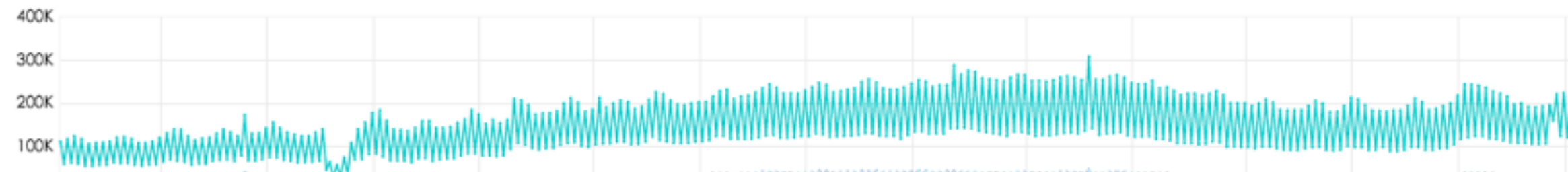
**3. 时间成本**

# 美团早期负载均衡结构

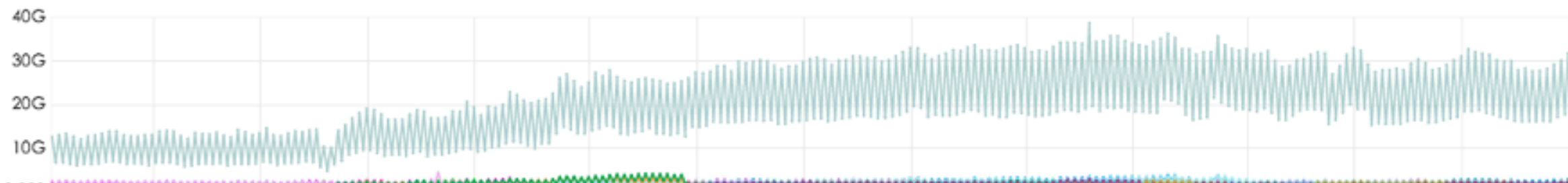
.....



# 美团流量增长情况



新建连接数 (cps) 100w -> 300w



吞吐量 (bps) 13G -> 38G

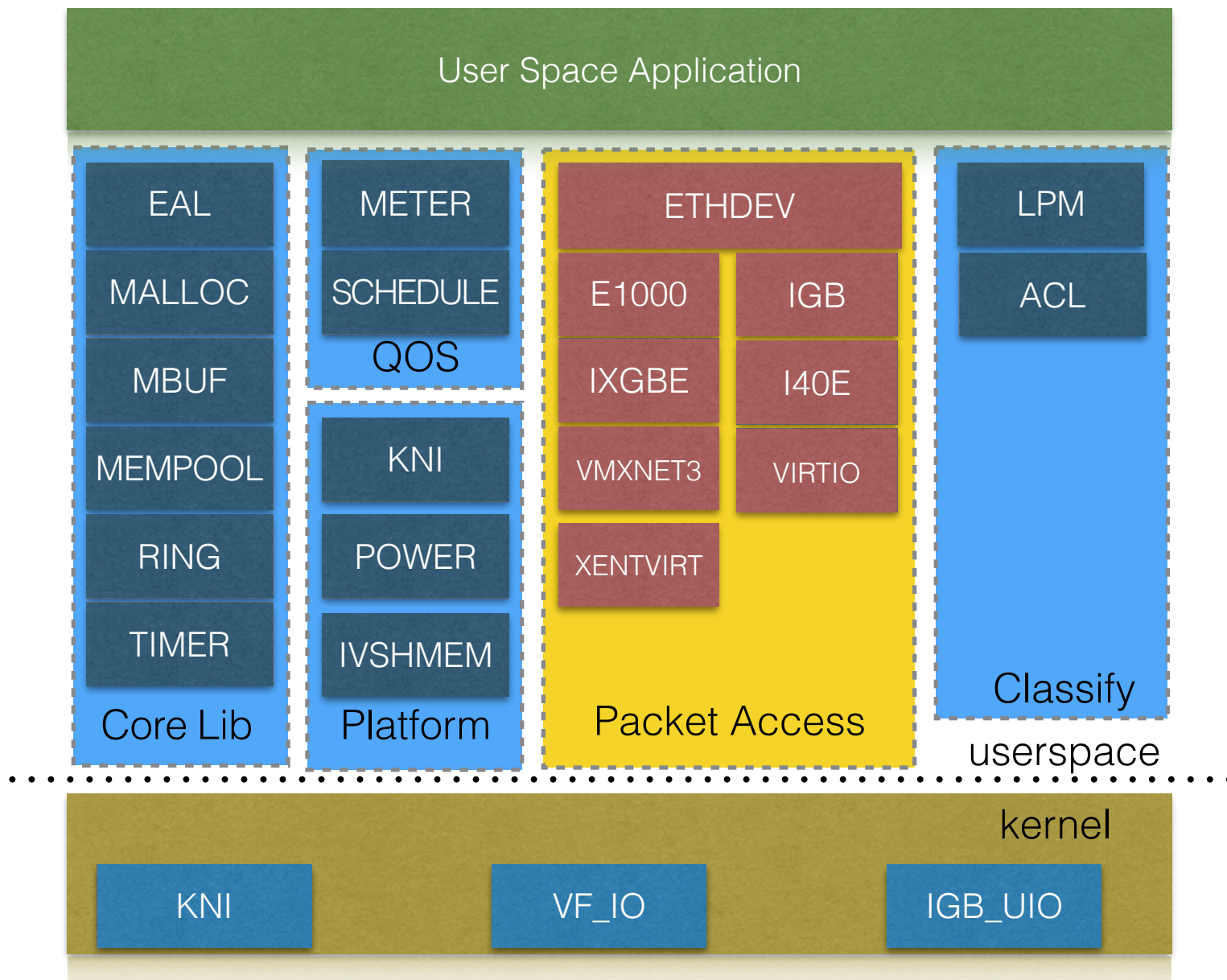
**！ 3倍的流量增长，LVS性能不足，故障率增加**

# 问题所在

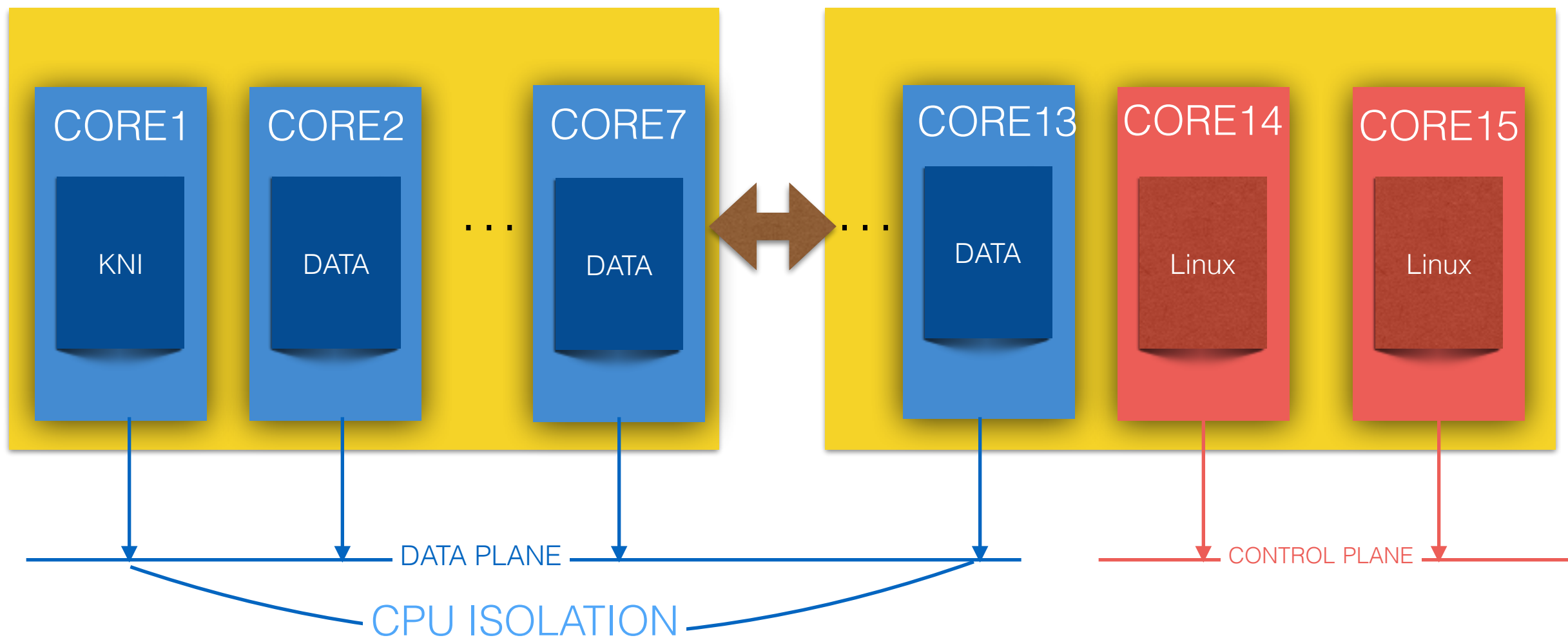
---

- |             |        |                  |
|-------------|--------|------------------|
| 1. 中断       | .....> | 1. PMD驱动         |
| 2. 过长的协议栈路径 | .....> | 2. kernel bypass |
| 3. 锁        | .....> | 3. 无锁的设计         |
| 4. 上下文切换    | .....> | 4. CPU绑定、隔离      |

# PMD驱动、kernel bypass -> DPDK

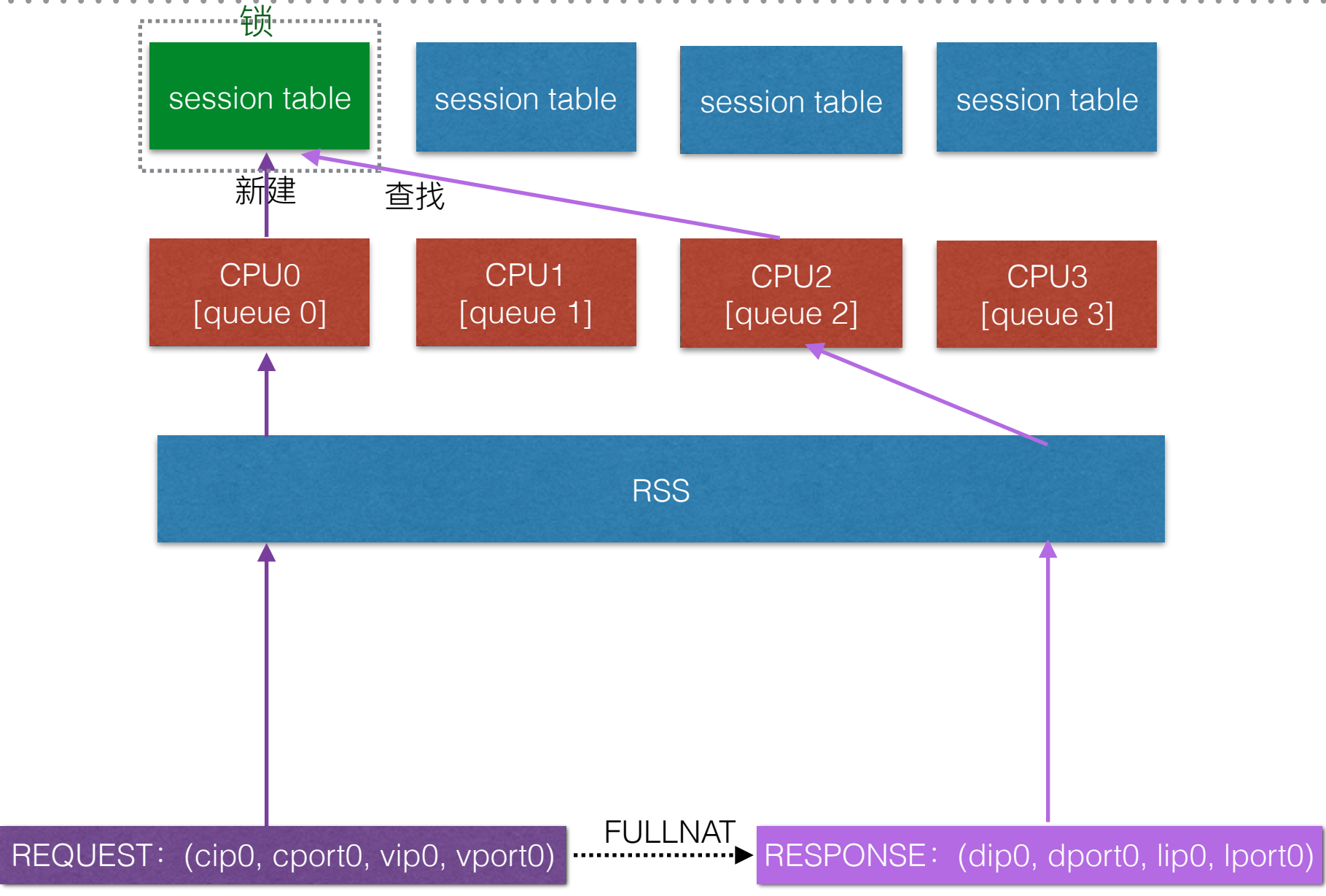


# 上下文切换

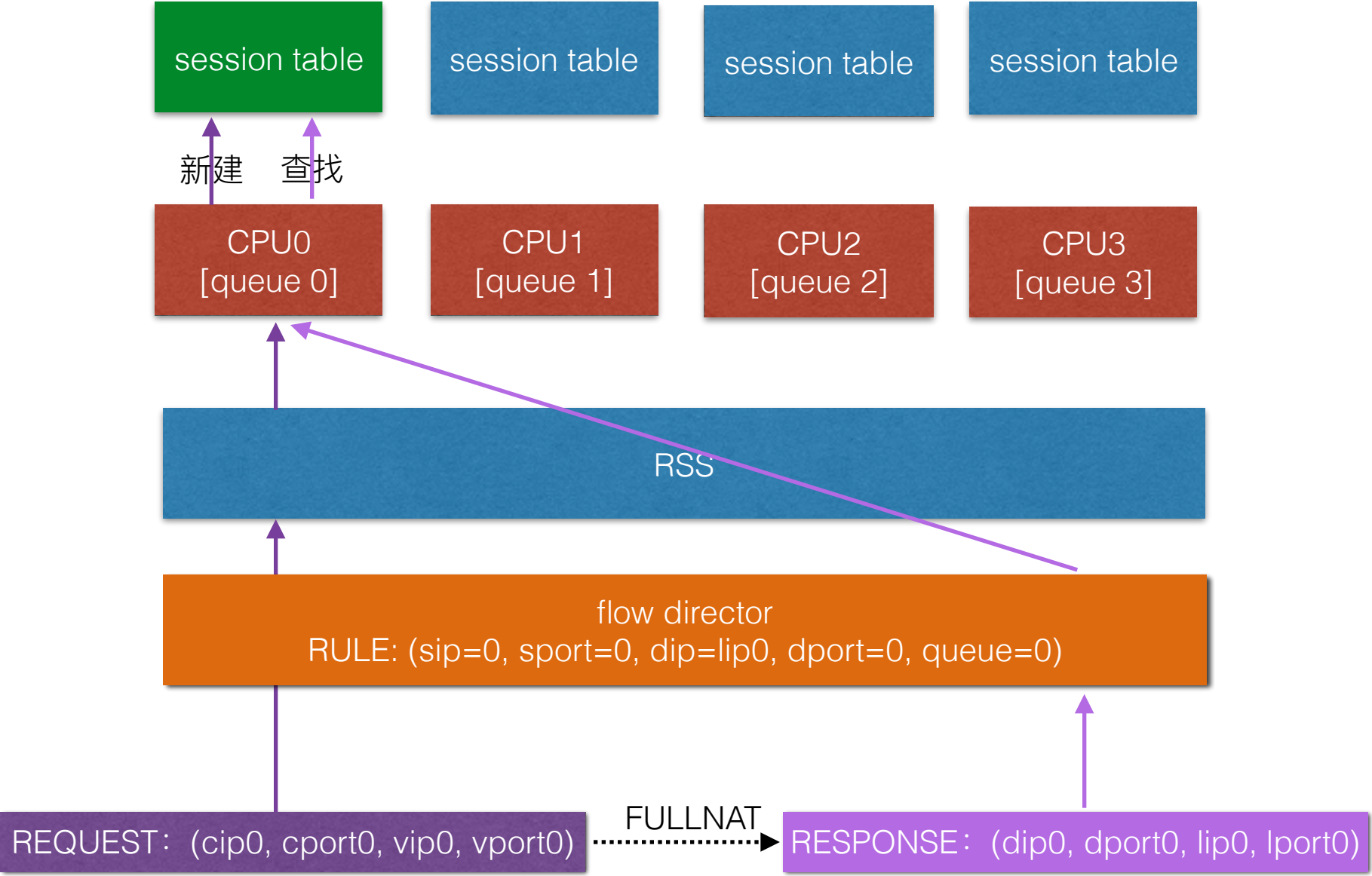




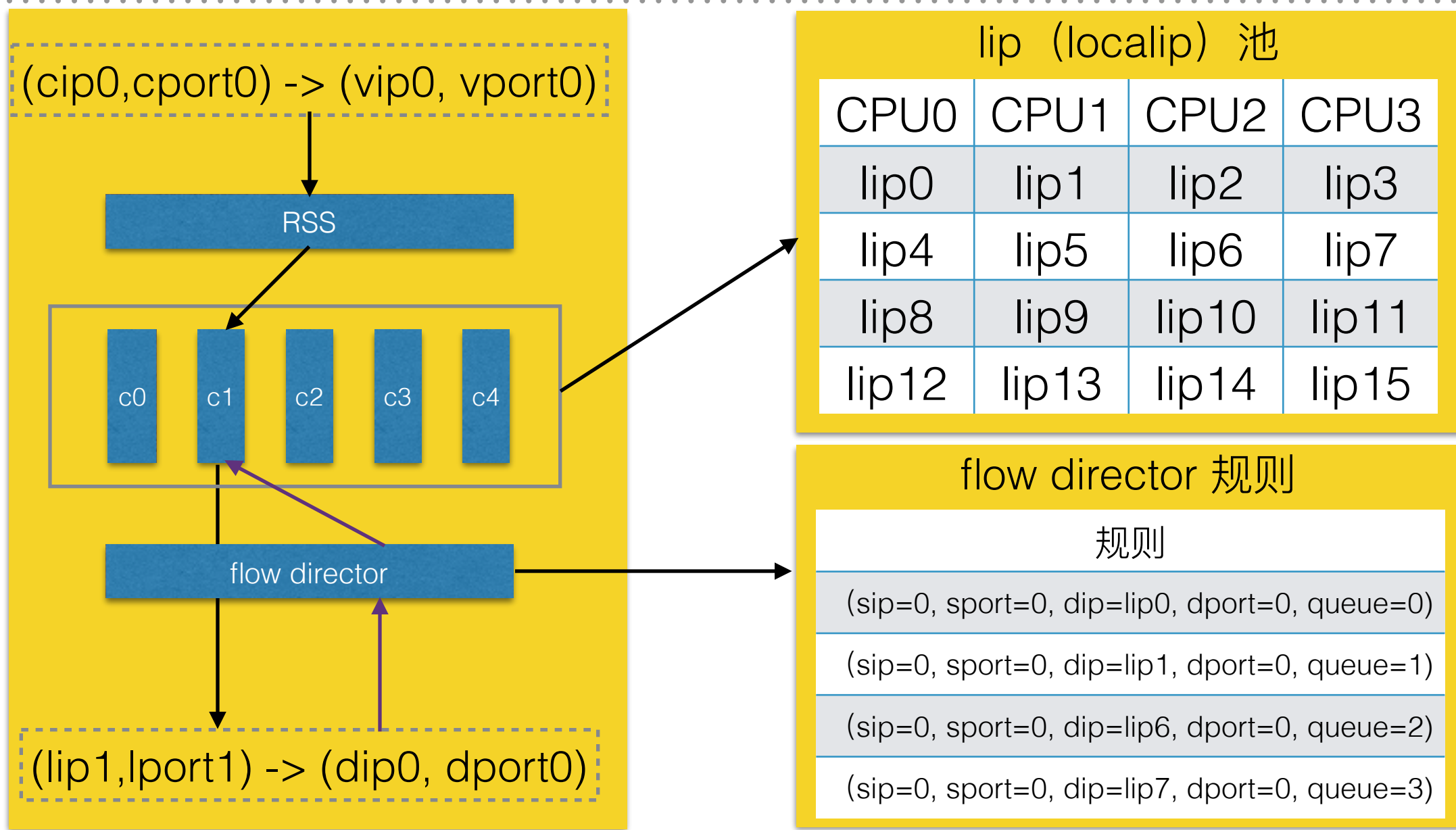
# 无锁的设计 - 地址转换的问题



# 无锁的设计 - 地址转换的问题



# 无锁的设计



# 性能测试

---

测试项	性能参数
<b>SYNPROXY</b>	2800w pps (67%CPU)
<b>http pps</b>	in 850w out 960w
<b>http bps (64bytes response)</b>	in 6G out 10G
<b>http qps</b>	380w

# 目录

---

- 负载均衡介绍
- 高性能
- 高可靠
- 技术展望

# 高可靠

.....

核心交换机  
(ECMP + OSPF)



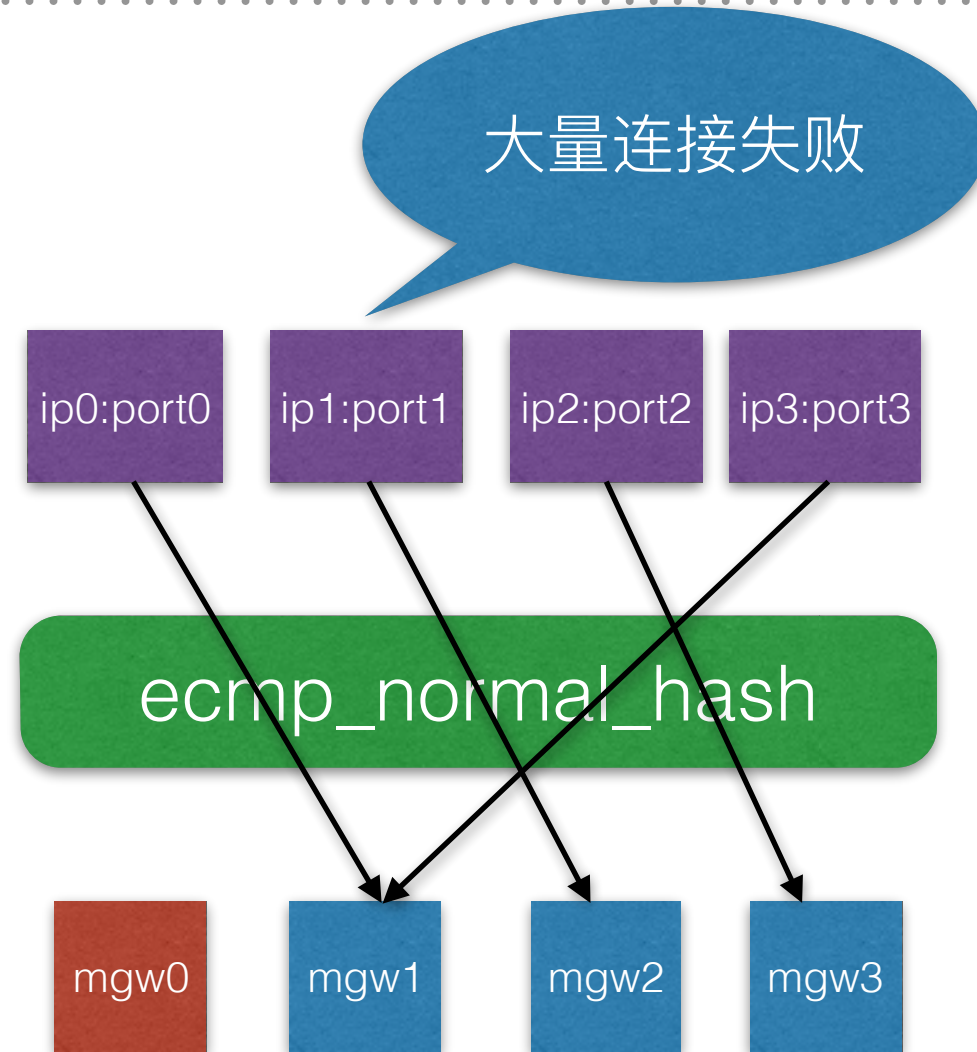
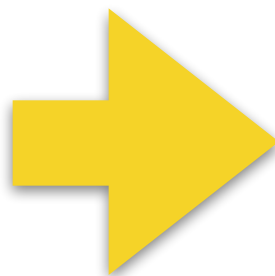
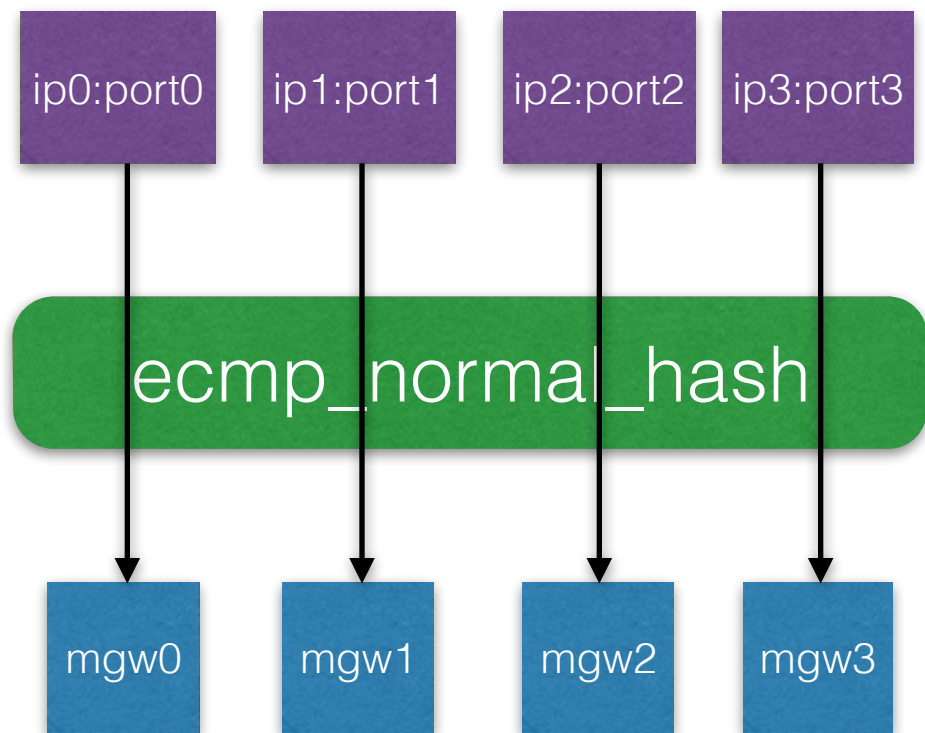
4层负载均衡  
(MGW)



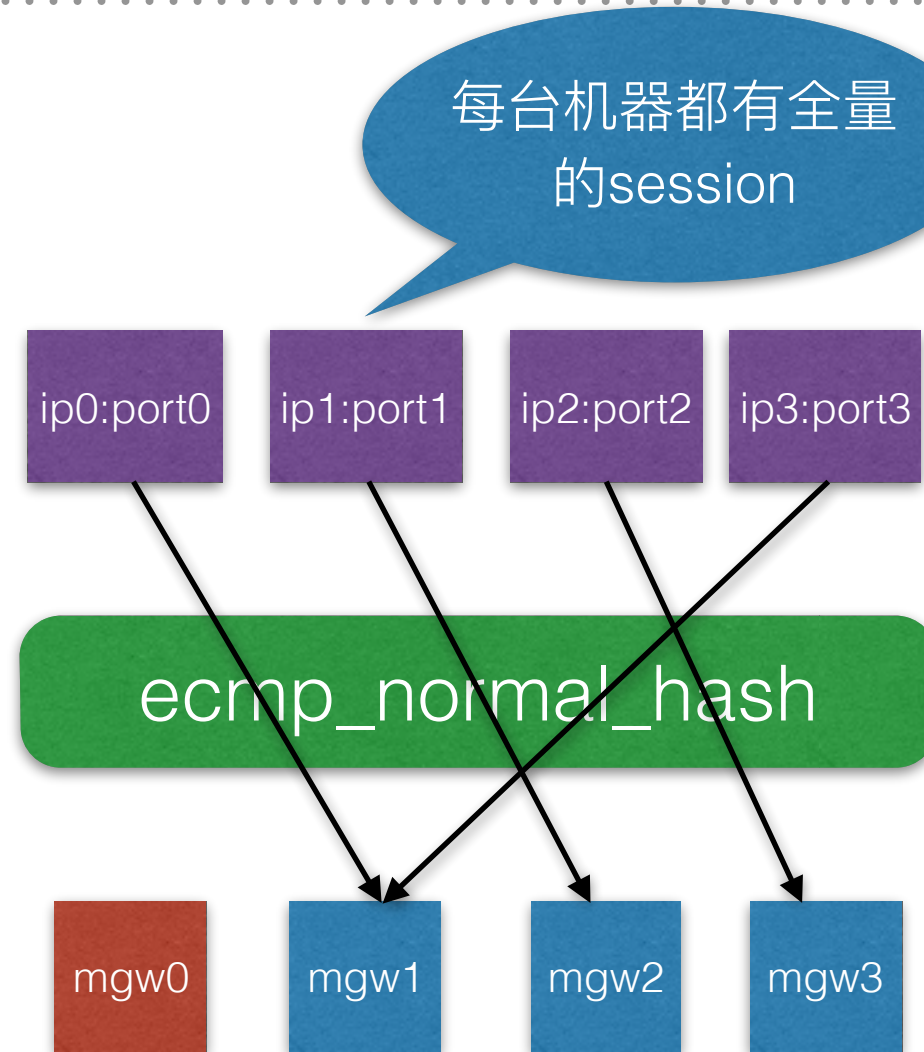
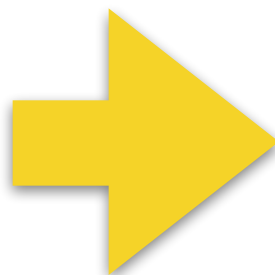
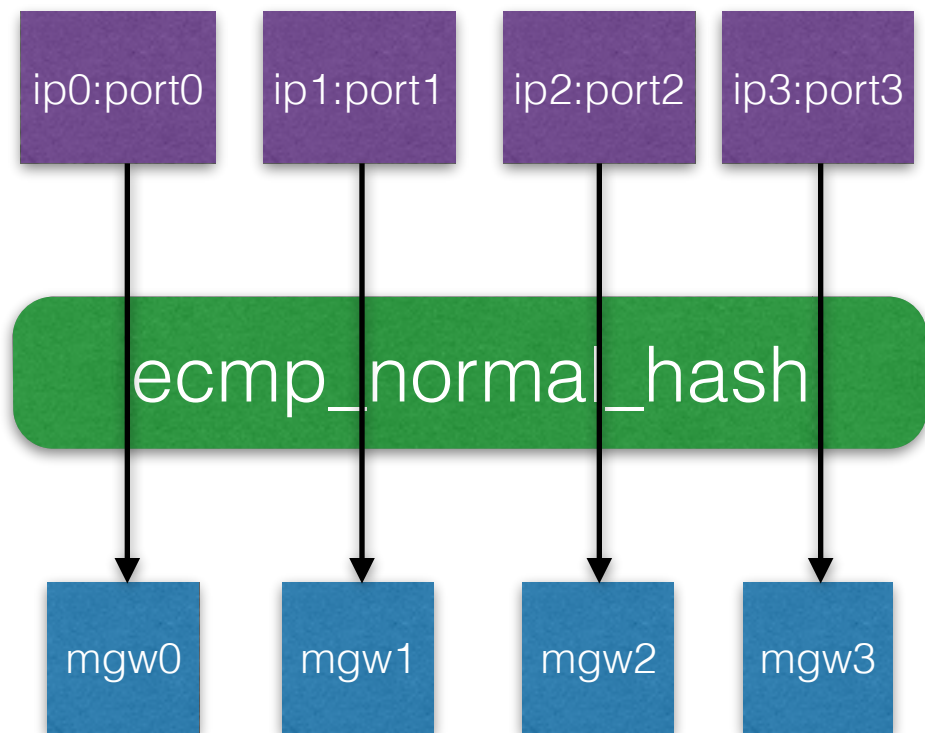
应用服务器



# 机器下线导致的问题



# session同步





# 故障检测与故障切换

## 故障切换

交换机侧不使用虚拟接口

半秒一次机器健康自检

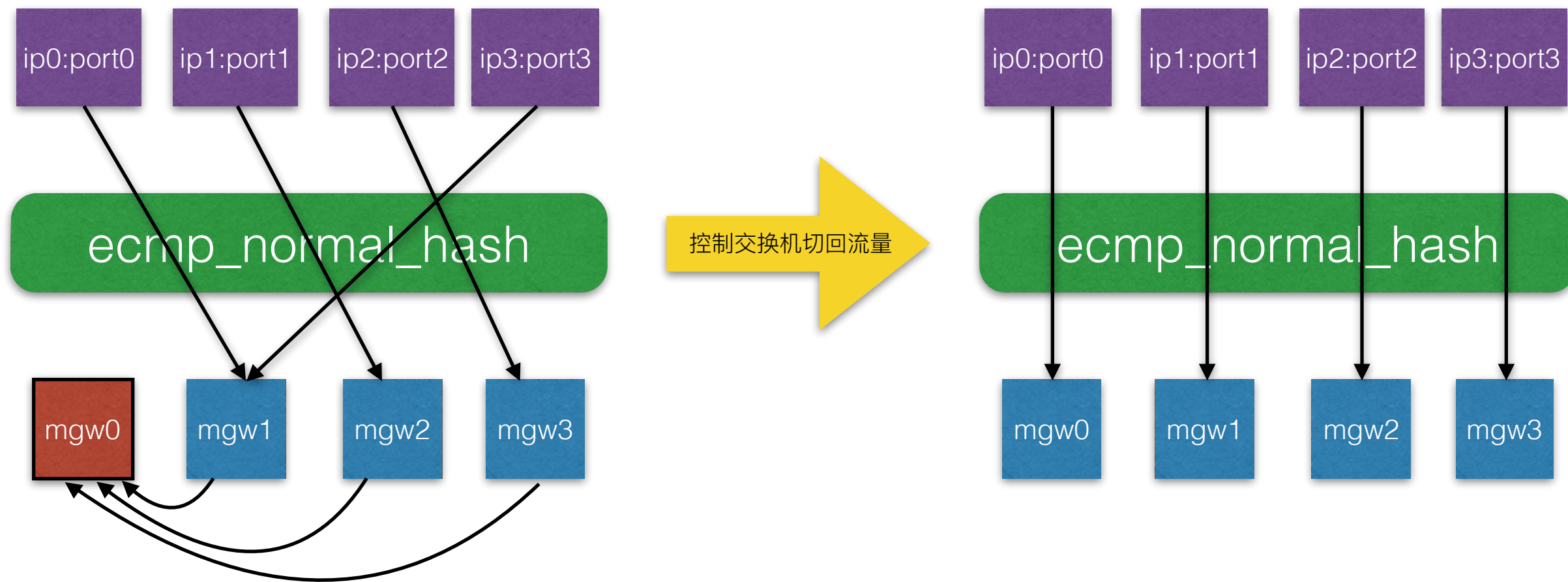
检测到故障自动给网口断电

捕获异常信号，物理网口断电

## 故障切换效果

测试程序发包间隔	100ms
升级操作丢包	0
主程序故障丢包	0
其他异常（网线等）	500ms

# 故障恢复与扩容



批量session同步旧session，同时通过增量同步接收新session

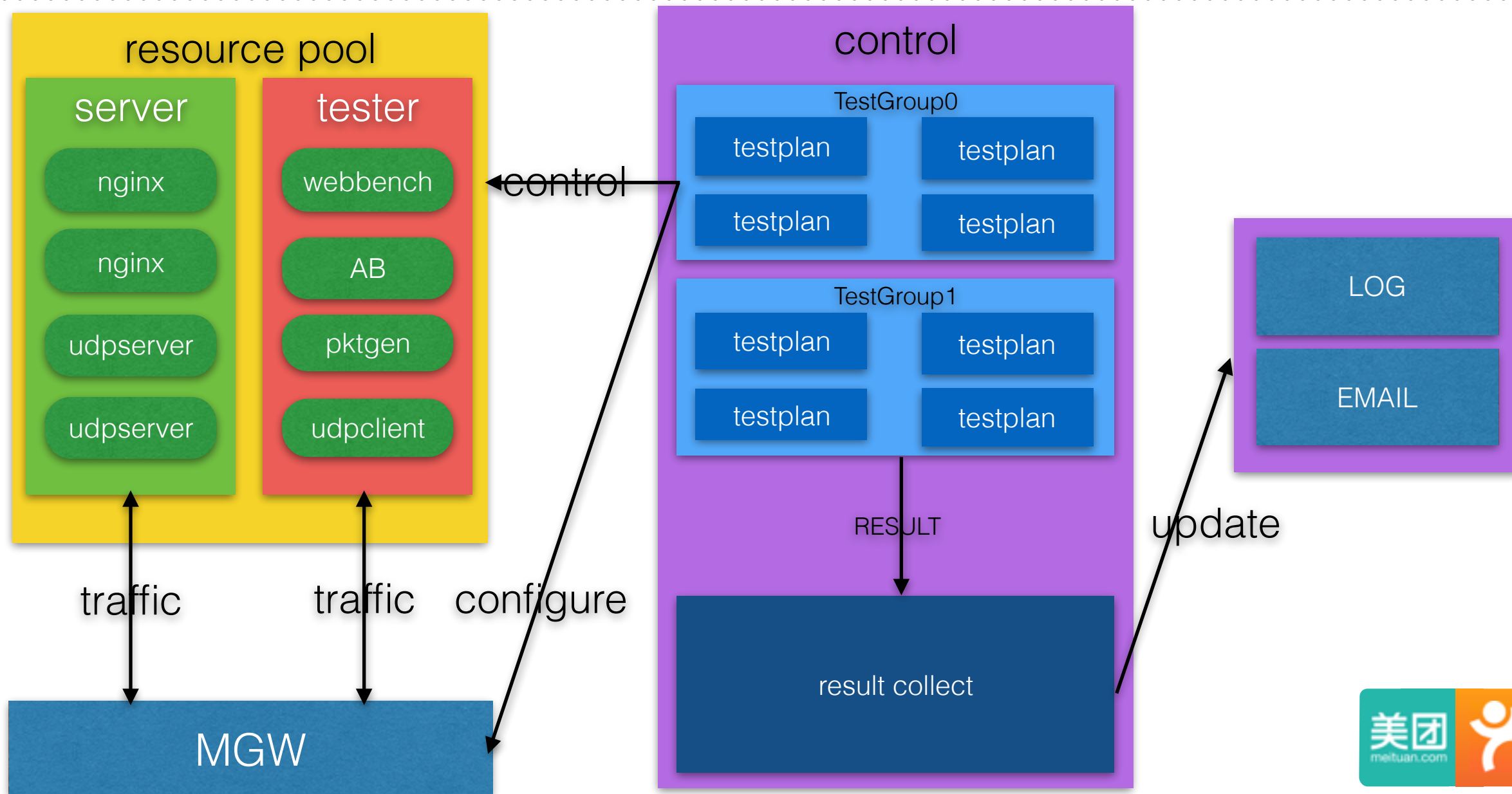
# MGW单机可靠性

---



自动化测试

# 自动化测试平台



# 自动化测试平台

---



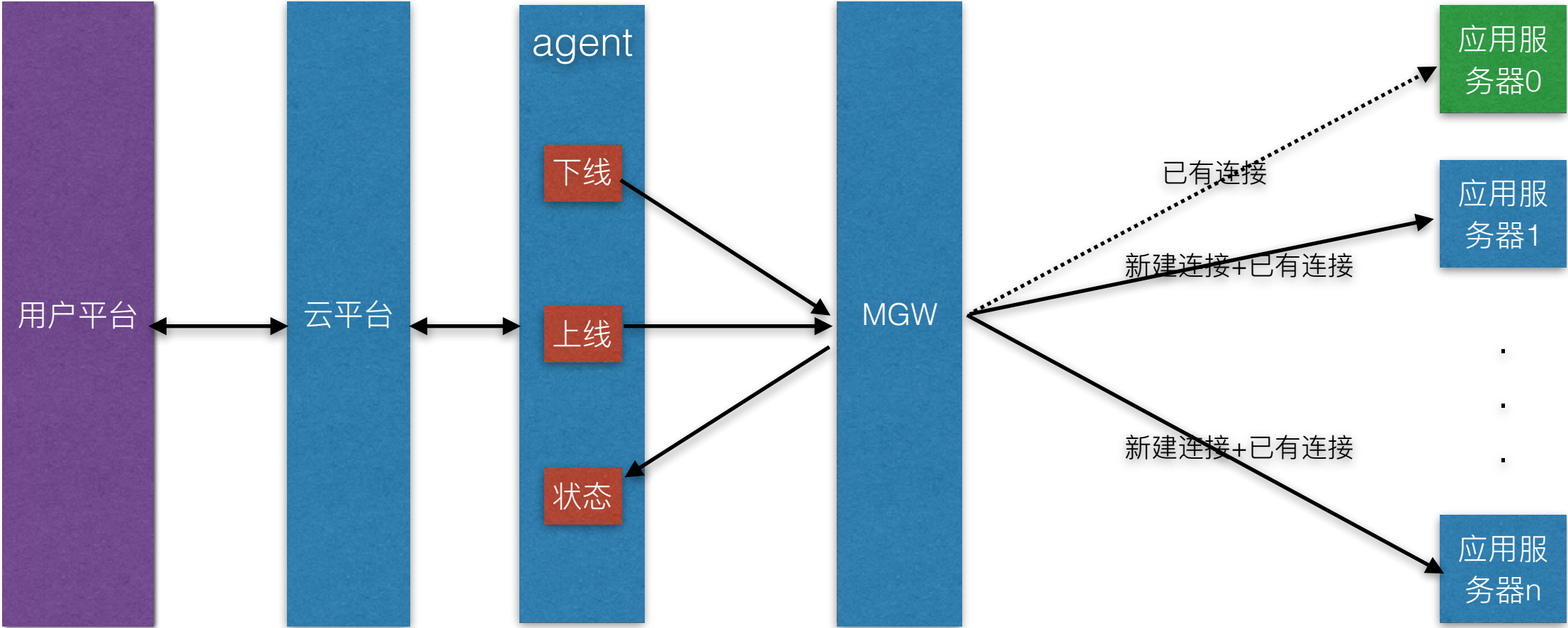
# 应用服务器可靠性

---



节点平滑下线

# 节点平滑下线



# 应用服务器可靠性

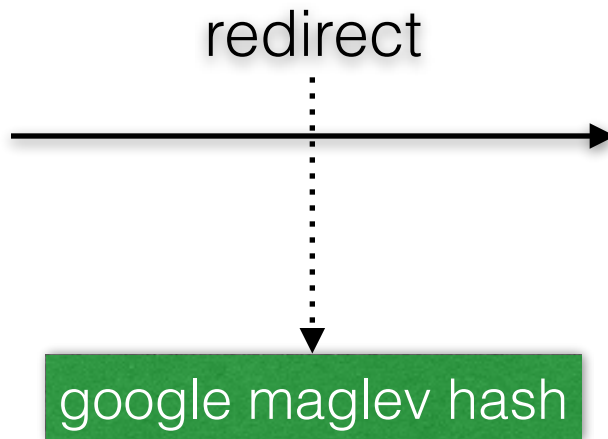
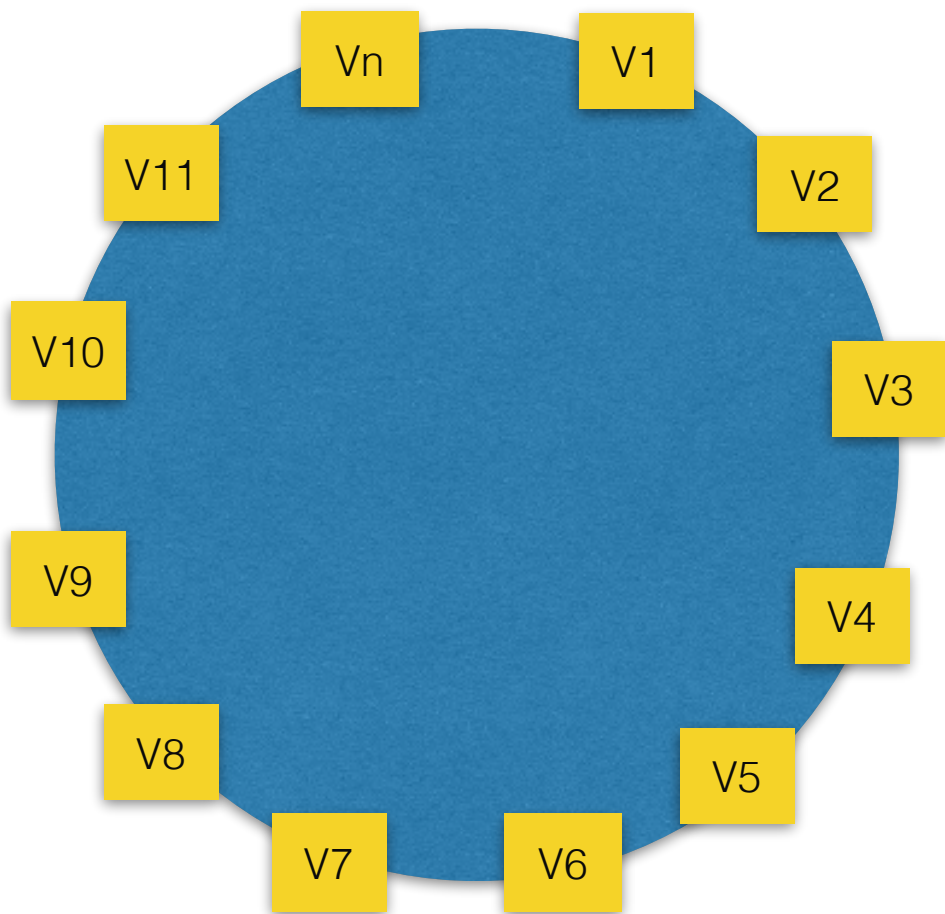
---



一致性源ip hash调度器



# 一致性hash调度器



V1	node0
V2	node1
V3	node2
V4	node3
V5	node0
V6	node1
V7	node2
V8	node3
V9	node0

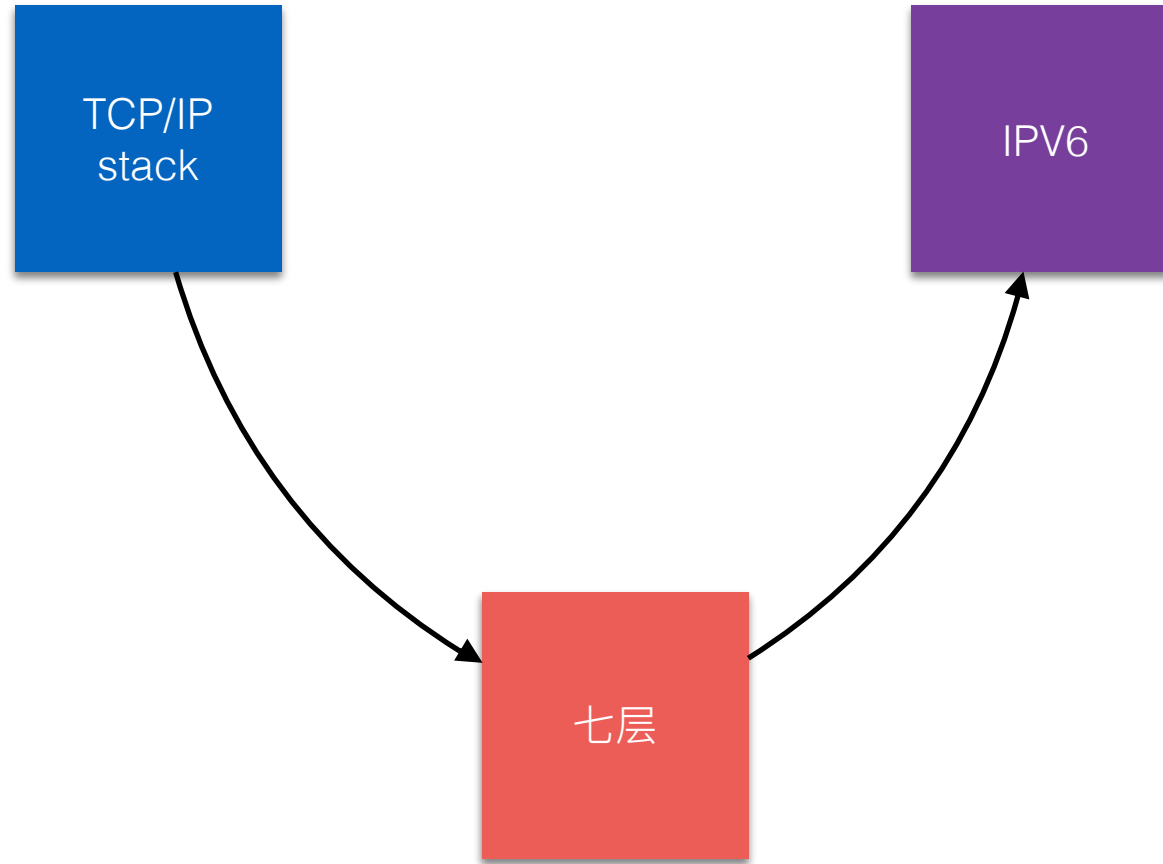
# 目录

---

- 负载均衡介绍
- 高性能
- 高可靠
- 技术展望

# 技术展望

---



谢谢大家

Q&A

